

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 08-241332

(43)Date of publication of application : 17.09.1996

(51)Int.Cl.

G06F 17/30

(21)Application number : 07-066727

(71)Applicant : FUJI XEROX CO LTD

(22)Date of filing : 02.03.1995

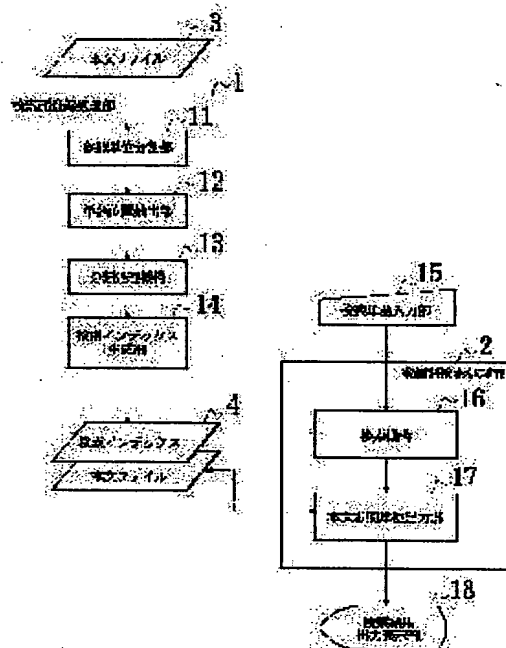
(72)Inventor : TATENO SHOICHI

(54) DEVICE AND METHOD FOR RETRIEVING ALL-SENTENCE REGISTERED WORD

(57)Abstract:

PURPOSE: To provide the device and method for retrieving all-sentence registered word which can obtain reference words of the main body of a tagged document by regarding tags as reference units for retrieval results and efficiently retrieving word positions in the main body.

CONSTITUTION: The device and method are equipped with a reference unit division part 11 which inputs a main body file 3 containing the main body of the tagged document and divides it into the reference units sectioned with the tags, a word position extraction part 12 which extracts pairs of words, included in the reference units, to be retrieved and the positions of the reference units in the main body where the words appear, a classification part 13 which classifies the pairs of the extracted words and the positions of the reference units according to the words and obtains a word position set of the positions of all the reference units where the words appear paired with the words, and a retrieval index generation part 14 which generates retrieval indexes for obtaining the position set from the words for the word position set.



BEST AVAILABLE COPY

LEGAL STATUS

[Date of request for examination] 18.10.1996

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number] 2896634

[Date of registration] 12.03.1999

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's

【特許請求の範囲】

【請求項 1】 タグを有する文書の本文を収めた本文ファイルを入力し、タグで区切られた参照単位に分割する参照単位分割部と、
参照単位に含まれる検索対象とする単語に対して、単語と当該単語が出現する本文における参照単位の位置の対を抽出する単語位置抽出部と、
抽出された単語と参照単位の位置の対を単語に従って分類し、単語に対し当該単語が出現する全ての参照単位の位置を組とした単語位置集合を得る分類部と、
単語位置集合に対し、単語から位置集合を得る検索インデックスを生成する検索インデックス生成部とを備えることを特徴とする全文登録語検索装置。

【請求項 2】 タグを有する文書の本文を収めた本文ファイルタグで区切られた参照単位に分割し、
参照単位内に含まれる検索対象とする単語に対して、単語と当該単語が出現する全ての参照単位の位置の対を抽出し、
抽出され単語と参照単位の位置の対を単語により分類し、
単語と当該単語が出現する全ての参照単位の位置を組とした単語位置集合を作成し、
作成された単語位置集合に基づいて、単語から位置集合を得ることができる検索インデックスを生成することを特徴とする全文登録語検索方法。

【請求項 3】 請求項 1 に記載の全文登録語検索装置において、更に、
検索インデックス生成部により作成された検索インデックスを用いて得られたタグの位置の集合に基づいて、参照単位を次のタグまであるいは適当な長さだけ表示して、検索結果を表示する検索処理部を有することを特徴とする全文登録語検索装置。

【請求項 4】 請求項 3 に記載の全文登録語検索装置において、
タグに付加されたフィールド情報に基づいて特定のフィールドを検索対象とする場合、
本文中の単語について、その単語の直前にフィールドの種別を示す文字列を付加したものを新たな単語とし、当該単語が出現する位置の直前にあるタグの位置との対を単語位置集合対とし、
単語を入力し、その語が出現する位置の直前にあるタグの位置の集合を結果として出力することを特徴とする全文登録語検索装置。

【請求項 5】 請求項 3 に記載の全文登録語検索装置において、
本文中に指定形式で記述された属性と値の対を含む場合、
前記属性と値の対を単語として登録し、
検索時に指定形式の一部を入力することにより、検索対象の文字列から始まる指定形式の属性と値の対を列挙す

2

ることを特徴とする全文登録語検索装置。

【請求項 6】 請求項 3 に記載の全文登録語検索装置において、
特定のフィールドを検索対象とし、本文中に指定形式で記述された属性と値の対を含む場合、
属性と値の対の直前にフィールドの種別を示す文字列を付加したものを単語として登録し、
検索時に指定形式の一部を入力することにより、検索対象の指定形式の文字列から始まる指定形式を列挙することを特徴とする全文登録語検索装置。

【発明の詳細な説明】

【0001】

【産業上の利用分野】 本発明は、全文登録語検索装置および方法に関し、特に、タグを有する文書において本文中の単語を登録し、全文の登録語の検索を能率よく行う全文登録語検索装置および方法に関するものである。

【0002】

【従来の技術】 従来から、ワークステーション上の文書編集装置（ワードプロセッサ）など、文書処理を行う文書編集装置においては、文書の作成を効率よく行うため、タグを用いて文書内容を部分的に区別して、予じめ、見出し、段落などの複数の文書部品を作成し、その各々の文書部品の間の関係を定めて、文書を構造化して編集することが試みられている。

【0003】 このような文書に対して構造の概念を取り入れた構造化文書の例としては、例えば、国際規格の ODA (ISO 8613: Open Document Architecture) や、SGML (ISO 8879: Standard Generalized Markup Language) の規格による構造化文書が知られている。ODA の規格による構造化文書を用いた文書処理方法の一例は、特開平 5-135054 号公報に記載されている「文書処理方法」が参照できる。

【0004】 ところで、SGML による構造化文書は、従来のテキスト処理システムとの親和性が高く、アメリカを中心として普及してきており、既に実用の段階に入っている。このような SGML による構造化文書の手法では、タグとよばれるマークを文書テキスト中に挿入することで、文書テキストを部分的に区分し（例えば文書部品として区分し）、その区分の関係を規定して構造化し、木構造の文書構造を表現する。

【0005】 SGML による構造化文書を例にとって、マーク（タグ）付けされた構造化文書の処理例について説明する。SGML による構造化文書では、予じめ文書の構造のひな型が与えられ、文書の構造は、その与えられたひな型の範囲内に制約される。この文書構造のひな型は、SGML においては、文書型定義 (DTD: Document Type Definition) と呼ばれる。

【0006】 SGML の構造化文書では、まず、文書型定義を規定して、文書の構造を表現するために、文書テキスト内にタグと呼ばれるマークを挿入し、そのタグに

(51) Int.Cl. ⁶	識別記号	序内整理番号	F I	技術表示箇所
G 0 6 F 17/30		9194-5L 9194-5L	G 0 6 F 15/401 15/40	3 1 0 A 3 7 0 A

審査請求 未請求 請求項の数6 F D (全 17 頁)

(21) 出願番号 特願平7-66727

(22) 出願日 平成7年(1995)3月2日

(71) 出願人 000005496

富士ゼロックス株式会社

東京都港区赤坂二丁目17番22号

(72) 発明者 館野 昌一

神奈川県横浜市保土ヶ谷区神戸町134番地

横浜ビジネスパークイーストタワー 富士ゼロックス株式会社内

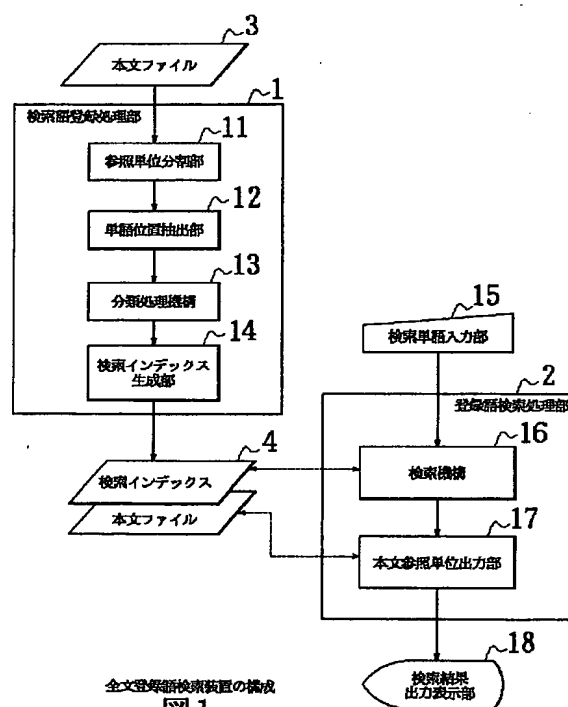
(74) 代理人 弁理士 南野 貞男 (外2名)

(54) 【発明の名称】 全文登録語検索装置および方法

(57) 【要約】

【目的】 タグ付き文書において、タグを検索結果の参照単位として、効率よく本文中の単語位置を検索し、本文の参照単語を得ることができる全文登録語検索装置および方法を提供する。

【構成】 タグを有する文書の本文を収めた本文ファイルを入力し、タグで区切られた参照単位に分割する参照単位分割部と、参照単位に含まれる検索対象とする単語に対して、単語と当該単語が出現する本文における参照単位の位置の対を抽出する単語位置抽出部と、抽出された単語と参照単位の位置の対を単語に従って分類し、単語に対し当該単語が出現する全ての参照単位の位置を組とした単語位置集合を得る分類部と、単語位置集合に対し、単語から位置集合を得る検索インデックスを生成する検索インデックス生成部とを備える。



より文書テキストを部分的に区分する。例えば、文書における一つの段落は、名前が“段落”とされたタグ<段落>を用いて、次のように表現される。

『<段落>これは一つの段落です。</段落>』

ここでのタグ<段落>が段落の開始を意味し、スタートタグと呼ばれる。タグ</段落>が段落の終了を意味し、エンドタグと呼ばれる。つまり、ここでは、タグの名前が“段落”とされたスタートタグ<段落>とエンドタグ</段落>との2つのタグを用いてマーク付けし、文書におけるテキストを文書部品として部分的に区分する。つまり、2つのタグの間に挟まれたテキスト部分が、タグで指示された構造の内容部分を示している。

【0007】名前が付けられたタグは各々が区別されて、文書型定義の中で構造上のその位置付けが定義される。その意味では、タグが文書の構造（構成要素）を表現している。したがって、混合が生じない場合において、以下で言う構造化文書（SGMLによる文書）の構造とは、タグと同義であることを意味している。

【0008】また、SGMLによる構造化文書（以下、SGML文書と略称する）においては、一部のタグを省略できる。その場合の省略の可／不可は、文書型定義（DTD）により指定する。省略はスタートタグおよびエンドタグのそれぞれに独立に指定できる。例えば、エンドタグ</段落>が省略可とする場合は、これが文書型定義内で指定された場合であり、その場合、先の例は『<段落>これは一つの段落です。』

と記述してもよいことになる。

【0009】SGML文書の文書型定義の具体例は、例えば、図10に示される。図10に示す文書型定義100により規定される文書構造では、名前が“題”とされたスタートタグ、“題”のエンドタグ、“段落”のエンドタグ、“図”のエンドタグ、および“図本体”のエンドタグが省略可能であることが定義されている。

【0010】更に、図10に示す文書型定義100の内容を具体的に説明すると、ここでの文書型定義（DTD）はSGMLの表記法に従って記述されているので、それに従って解釈できる。つまり、文書型定義の内容の行の最初の“<!”はマークアップ宣言区切り子であり、空白なしに続く次の“ELEMENT”は要素宣言キーワードである。この行の最初の“<!ELEMENT”により、次に続く記述によって、その構造の内容（下部の構造）がどのようなものかを指定する。そして、その次に記述される項目の名前（文書、章、題、段落、図など）が、対象となるタグの名前を表している。

【0011】更に、次の記号（“-”，“-”，“0”，“0 0”など）は、その項目の対象のタグが、スタートタグおよびエンドタグの順でそれぞれ省略可能かどうかを表す記号である。“-”が省略不可を意味し、“0”が省略可を意味する。例えば、ここでの記号が“- 0”であれば、スタートタグは省略不可であ

り、エンドタグは省略可であることを意味する。

【0012】更に続く次の項目は、タグの下部の構造を表す規定の定義である。ここでの記号“|”は項目（タグ）が順序立てて出現することを意味し、記号“|”はどちらかの項目であれば良いことを意味し、記号“*”は0回以上の繰り返しを意味する。また、記号“?”はそれがあってもなくても良いことを意味している。

【0013】したがって、例えば、タグの下部構造が“（章題，段落*，章*）”と規定されている場合は『章題の次に段落の0回以上の繰り返しがあり、更にその次に章の0回以上の繰り返しがある』という順序で下部の構造が規定されることを意味する。具体例で説明すると、図10に示す文書型定義100の第2行目のように、タグの下部構造が“（題，（段落|図）*，章*）”と規定されている場合、『章の次に段落または図の0回以上の繰り返しがあり、その次に章の0回以上の繰り返しがある』順序で下部の構造が規定されることを意味する。

【0014】また、第3行目および第4行目に記述されているタグの下部構造の“#PCDATA”はSGMLの予約語の1つであり、構造の規定で、その内容が文字データであること意味している。したがって、文書型定義100の例では、章を構成する「題」と「段落」のタグの下部には文字データが来ることを意味している。

【0015】つまり、図10に示す文書型定義（DTD）の意味するところによる文書構造のひな型では、当該文書が、「章」の繰り返しからなる“<文書>”というタグから始まる文書であり、その「章」は、「題」の次に「段落」または「図」の0回以上の繰り返しがあり、更にその次に「章」の0回以上の繰り返しがあるものから構成されている。そして、ここでの「題」および「段落」が、文字データから構成される。

【0016】更に細部の規定として、文書構造の「図」の内容は、「題」とそれに続く「図本体」から構成されると定義され、「図本体」は、例えば外部のイメージファイルを参照するので、下部構造を持たない（“EMPTY”）と定義される。また、ここでの構造のタグの省略可／省略不可の指定は、「文書」および「章」はタグの省略不可であり、「段落」，「図」，および「図本体」は、エンドタグのみが省略可であり、「題」は両方のタグが省略可であるということが定義される。

【0017】このような文書型定義に従っている実際の文書（以下、対象文書と呼ぶ）の例としては、例えば、図11に示されているSGML文書110がある。なお、この図11に示す文書の中では、文書の構造の深さに応じてインデントーションを変えて表記しているが、この表記は、ここでの構造化文書の文書例の説明上、見やすくするために行っているものであり、実際の文書ではインデントーションされないことが多い。

【0018】図11を参照すると、この例の構造化文書

のSGML文書110を見ると理解されるように、
「章」を構成するその下位の構造の「題」を表すタグは、スタートタグおよびエンドタグは共にこの文書中には現れていない。しかし、実体的には、第2行目のタグ“<章>”とその内容部分の“SGMLとは”との間のスタートタグ“<題>”が省略されている形となっている。なお、このようなタグが省略されているかどうかは、前述の文書型定義100を参照しなければ分からない。したがって、対象文書は常にそれに従っている文書型定義とのセットでないと正確な構造を読み取ることができない。

【0019】SGML文書では、このようにタグの省略が行われるため、SGML文書进行处理するには、まず、文書構造を解析する処理（SGMLパーサによる構文解析処理）が必要とされる。文書構造の解析の処理では、対象文書を解析しながら、文書型定義との照らし合わせを行い、対象文書において省略されたタグを復元する処理が主として行われる。実際の文書処理において実行される構文解析処理では、これ以外の処理（例えば、属性の復元やエンティティの展開などの処理）も行われる。

【0020】図11に例示したSGML文書110を対象文書として、タグ（構造）の復元処理を行うと、図12に示すようなSGML文書120が得られる。図12に示すSGML文書120においては、下線を引いた部分が復元されたタグ（構造）を示している。この対象文書は、図10に示すような文書型定義100を参照して、省略されたタグが復元されたものである。つまり、「章」の構造の規定から、タグ<章>の次には必ずタグ<題>が存在しなければならないので、まず、タグ<章>の次にタグ<題>を復元し、同様に、タグ<図>の次には必ずタグ<題>が存在しなければならないので、タグ<図>の次にタグ<題>を復元する。また、それぞれのエンドタグが省略されているので、内容部分の次に（対応する次のタグの手前の位置に）、それぞれのエンドタグ</題>、</段落>などを復元する。このようにして下線が引かれたようにタグ（構造）が復元される。

【0021】次に、このようにタグが復元され、構造表現されたSGML文書120において、構造を検索する場合の処理について説明する。構造化文書においては、文書編集を行う場合、単なるテキストの文字列の検索だけでなく、文書構造における構造の検索も文書処理の重要な処理となる。これは、構造化文書の処理を行う場合においては、文書構造の構造を利用した編集処理が積極的に行なわれるためである。

【0022】構造化文書の検索においては、従来のような文字列の検索だけではなく、構造を積極的に利用した検索も有効に利用される。例えば、文書内のSGMLに関連した図を検索したい場合、これまでの検索処理で

は、文書内を全文検索（文字列検索）を行い、テキストの文字列からその「関連した図」の文字列を捜し出していた。

【0023】しかし、文書構造の構造自体を検索に用いると、例えば、「図のタイトルにSGMLを含んでいる図」や「図の下部構造の題」のように文書構造における構造を指示して、検索を行うことができ、より対象を絞った検索を行うことができることになる。また、その場合の検索処理も、文書の構造に従って対象範囲が限定して検索できるので、検索処理の効率が良くなるという利点も持つ。

【0024】前述したように、SGML文書は、テキスト中にマーク付けを行うだけのタグを埋め込んだ形式の文書アーキテクチャとなっているため、従来からのテキスト処理システムとの親和性が高く、構造がマーク付けのタグで表現されるため、構造を検索する際にも特別な装置あるいは処理プログラムを用いなくとも良い。タグの文字列を検索するという文字列検索を用いて、文書構造の検索を行うことができる。つまり、従来からのテキスト処理装置（文書エディタなど）により、SGML文書を作成することができ、基本的にはスタートタグとそれに対応するエンドタグを、タグの文字列を検索するという従来の文字列検索のテキスト検索手法を用いて検索することにより、構造検索を行うことができる。

【0025】

【発明が解決しようとする課題】このように、SGMLなどのタグ付き文書の規格が標準化され、多方面で利用され始めている。このタグ付き文書は、フォーマット情報や、その他の文書に関する構造の情報を示すマークアップをタグにより表現する手法を取り入れた文書であるため、SGMLの標準化の規格によって、文書の内容が構造化されて、その内容の利用が容易になる。このため、企業、官庁、学校の内外を問わず、この種の文書の流通が盛んになり始めている。このように、タグ付き文書が電子化されて、蓄積されていくことにより、大規模な電子化文書の資源が蓄えられ、利用可能になる。

【0026】しかし、SGMLなどのタグ付き文書の中から、必要な情報を探する場合、文書構造は、タグを利用することにより容易に判定されるので、構造の検索は容易になっているが、文書内容については、これまでのフルテキストサーチなどの手法を利用しなければならず、十分に効率よく内容の検索までは行えないという問題があった。なお、タグを利用することにより、文書構造の位置関係など判別できるので、このようなタグを利用して、タグを検索結果の参照単位の区切りとすることができ、文書内容の利用が効率化できる。

【0027】本発明は、上述のような問題を解決するためになされたものであり、本発明の目的は、タグ付き文書において、タグを検索結果の参照単位として、効率よく本文中の単語位置を検索し、本文の参照単語を得るこ

とができる全文登録語検索装置および方法を提供することにある。

【0028】

【課題を解決するための手段】上記のような目的を達成するため、本発明の第1の特徴とする全文登録語検索装置は、タグを有する文書の本文を収めた本文ファイルを入力し、タグで区切られた参照単位に分割する参照単位分割部と、参照単位に含まれる検索対象とする単語に対して、単語と当該単語が出現する本文における参照単位の位置の対を抽出する単語位置抽出部と、抽出された単語と参照単位の位置の対を単語に従って分類し、単語に
10 対し当該単語が出現する全ての参照単位の位置を組とした単語位置集合を得る分類部と、単語位置集合に対し、単語から位置集合を得る検索インデックスを生成する検索インデックス生成部とを備えることを特徴とする。

【0029】また、本発明において、第2の特徴とする全文登録語検索方法は、タグを有する文書の本文を収めた本文ファイルをタグで区切られた参照単位に分割し、参照単位内に含まれる検索対象とする単語に対して、単語と当該単語が出現する全ての参照単位の位置の対を抽出し、抽出され単語と参照単位の位置の対を単語により
20 分類し、単語と当該単語が出現する全ての参照単位の位置を組とした単語位置集合を作成し、作成された単語位置集合に基づいて、単語から位置集合を得ることができる検索インデックスを生成することを特徴とする。

【0030】また、本発明の第3の特徴とする全文登録語検索装置では、更に、検索インデックス生成部により作成された検索インデックスを用いて得られたタグの位置の集合に基づいて、参照単位を次のタグまであるいは
30 適当な長さだけ表示して、検索結果を表示する検索処理部を有することを特徴とする。

【0031】また、本発明の第4の特徴とする全文登録語検索装置は、タグに付加されたフィールド情報に基づいて特定のフィールドを検索対象とする場合、本文中の単語について、その単語の直前にフィールドの種別を示す文字列を付加したものを新たな単語とし、当該単語が出現する位置の直前にあるタグの位置との対を単語位置集合対とし、単語を入力し、その語が出現する位置の直前にあるタグの位置の集合を結果として出力することを
40 特徴とする。

【0032】また、本発明の第5の特徴とする全文登録語検索装置は、本文中に指定形式で記述された属性と属性の値の対を含む場合、属性と値の対を単語として登録し、検索時に指定形式の一部を入力することにより、検索対象の指定形式の文字列から始まる指定形式の属性と値との対を列挙することを特徴とする。

【0033】また、本発明の第6の特徴とする全文登録語検索装置は、特定のフィールドを検索対象とし、本文中に指定形式で記述された属性と値の対を含む場合、属性と値の対の直前にフィールドの種別を示す文字列を付
50

加したものを単語として登録し、検索時に指定形式の一部を入力することにより、検索対象の指定形式の文字列から始まる指定形式を列挙することを特徴とする。

【0034】

【作用】本発明の第1の特徴とする全文登録語検索装置においては、タグを有する文書に対して、まず、参照単位分割部が、タグを有する文書の本文を収めた本文ファイルを入力し、タグで区切られた参照単位に分割すると、単語位置抽出部が、参照単位に含まれる検索対象とする単語に対して、単語と当該単語が出現する本文における参照単位の位置の対を抽出する。次に、分類部が、抽出された単語と参照単位の位置の対を単語に従って分類し、単語に
10 対し当該単語が出現する全ての参照単位の位置を組とした単語位置集合を得る。そして、検索インデックス生成部が、単語位置集合に対し、単語から位置集合を得る検索インデックスを生成する。これにより、文書内の全ての単語は、その直前のタグの位置と共に、検索インデックスに登録されるので、検索インデックスを用いて検索を行うことによって、文書中の単語の参照単位のタグ位置が直ちに検索でき、高速に参照単位の本文部分が表示出力されることになる。

【0035】また、本発明の第2の特徴とする全文登録語検索方法においては、検索対象の文書に対して、検索処理を開始する前に、まず、タグを有する文書の本文を収めた本文ファイルをタグで区切られた参照単位に分割し、参照単位内に含まれる検索対象とする単語に対して、単語と当該単語が出現する全ての参照単位の位置の対を抽出し、抽出され単語と参照単位の位置の対を単語により分類し、単語と当該単語が出現する全ての参照単位の位置を組とした単語位置集合を作成し、作成された
30 単語位置集合に基づいて、単語から位置集合を得ることができる検索インデックスを生成する。これにより、検索インデックスが作成された後は、次の検索処理から、作成された検索インデックスを利用することにより、直ちに検索対象の単語の参照単位のタグ位置が直ちに検索でき、高速に参照単位の文書部分が表示出力されることになる。

【0036】また、本発明の第3の特徴とする全文登録語検索装置においては、検索処理部が、検索インデックス生成部により作成された検索インデックスを用いて得られたタグの位置の集合に基づいて、参照単位を次のタグまであるいは適当な長さだけ表示して、検索結果を表示する。すなわち、検索対象の単語を検索する場合、検索インデックス生成部により作成された検索インデックスを利用することにより、検索結果として、単語に対してタグの位置の集合が得られるので、これにより、得られたタグの位置の集合に基づいて、参照単位を次のタグまであるいは適当な長さだけ表示して、検索結果を表示する。このため、検索結果の対象文書の箇所の表示が能率よく行える。
50

【0037】また、本発明の第4の特徴とする全文登録語検索装置は、タグに付加されたフィールド情報に基づいて特定のフィールドを検索対象とする場合、本文中の単語について、その単語の直前にフィールドの種別を示す文字列を付加したものを新たな単語とし、当該単語が出現する位置の直前にあるタグの位置との対を単語位置集合対とする。そして、検索する単語を入力し、その単語が出現する位置の直前にあるタグの位置の集合を結果として出力する。これにより、タグに付加されたフィールド情報に基づいて特定のフィールドを検索対象とすることができる。

【0038】また、本発明の第5の特徴とする全文登録語検索装置においては、更に、本文中に指定形式で記述された属性と属性の値の対を含む場合、属性と値の対を単語として登録し、検索時に指定形式の一部を入力することにより、検索対象の指定形式の文字列から始まる指定形式の属性と値の対を列挙する。これによって、本文の中において、書式など指定形式で記述された属性と属性の値の対を含む場合においても、それを検索対象とした検索が可能となる。

【0039】本発明の第6の特徴とする全文登録語検索装置においては、特定のフィールドを検索対象とし、本文中に指定形式で記述された属性と値の対を含む場合、属性と値の対の直前にフィールドの種別を示す文字列を付加したものを単語として登録し、検索時に指定形式の一部を入力する。これによって、検索対象の指定形式の文字列から始まる指定形式を列挙する。このため、特定のフィールドを検索対象とし、本文中に指定形式で記述された属性と値の検索に対しても容易に対応できる。

【0040】このようにして、本発明の全文登録語検索装置によれば、本文中の検索対象とする単語について、例えば、本文中の全ての単語について、その単語が出現する位置の直前にあるタグの位置を全て集めて、検索インデックスを作成する。この検索インデックスを用いることにより、検索処理を行う場合、検索対象の単語を入力とし、その単語が出現する位置の直前にあるタグの位置の集合が、検索結果として出力できる。つまり、タグを有する文書において、タグを検索結果の参照単位の区切りとすると、文書内に表われる全ての単語が、その直前のタグの位置と共に、検索インデックスとして保存される。このため、文書内の全ての単語に対し、当該単語を含むタグで区切られた参照単位を即座に検索することが可能となる。

【0041】

【実施例】以下、本発明の一実施例を図面を用いて具体的に説明する。図1は本発明の実施例の全文登録語検索装置の装置構成の要部を示すブロック図である。図1において、1は検索語登録処理部、2は登録語検索処理部、3はタグを有する文書の本文ファイル、4は本文ファイルに付加された検索インデックス、11は参照単位

分割部、12は単語位置抽出部、13は分類処理機構、14は検索インデックス生成部、15は検索単語入力部、16は検索機構、17は本文参照単位出力部、18は検索結果出力表示部である。

【0042】ここでの全文登録語検索装置においては、検索対象の単語を入力して本文検索の処理を実行する前に、その前処理として、検索語登録処理部1が、タグを有する文書の本文ファイル3から、検索対象とする単語の登録を行い、本文ファイルに付加する検索インデックス4を作成する。検索インデックス4が付加された本文ファイルは、登録語検索処理部2において、本文ファイルに付加された検索インデックス4を利用して、検索対象の本文からその登録語を検索する処理が行われる。図1を参照して説明する。

【0043】検索語登録処理部1において、まず、参照単位分割部11にタグを有する文書の本文ファイル3を入力する。参照単位分割部11は、本文ファイル3が入力されると、本文をタグで区切られた参照単位に分割する。この参照単位の本文を入力として、次に、単語位置抽出部12が、検索対象とする単語として、単語と当該単語が出現する本文における参照単位の位置の対を抽出する。

【0044】次に、分類処理機構13が、抽出された単語と参照単位の位置の対を単語に従って分類し、後述するように、単語に対し当該単語が出現する全ての参照単位の位置を組とした単語位置集合を生成する。そして、検索インデックス生成部14が、得られた単語位置集合に対し、単語から位置集合を得る検索インデックス4を生成し、本文ファイル3に対して、その対応する検索インデックス4を付加して、検索インデックス4付きの本文ファイルを作成する。これにより、ここでの検索インデックス4が付加された本文ファイルは、検索インデックス4を用いる登録語検索処理部2の検索処理により、検索対象の単語から高速に本文ファイルのタグで区切られた参照単位の位置を得ることができ、該当の参照単位の内容を直ちに表示できる。

【0045】検索対象の単語から本文ファイルの検索を行う場合、登録語検索処理部2においては、検索単語入力部15を介して、検索対象とする単語を入力すると、検索機構16が、本文ファイルに付加された検索インデックス4を用いて、検索対象の単語の検索処理を行い、その単語の対応の参照単位の位置の集合を検索する。参照単位の位置が検索できると、次に、本文参照単位出力部17が、その参照単位の位置から本文ファイルをアクセスして、該当の参照単位を直ちに出力し、検索結果出力表示部18を介して、該当の参照単位を出力表示する。

【0046】このように、タグを有する文書の本文ファイル3を、検索語登録処理部1の参照単位分割部11に入力すると、参照単位分割部11は、タグで区切られた

参照単位に分割し、参照単位分割部11から参照単位とその位置を得る。単語位置抽出部12は、この参照単位を入力とし、この参照単位内に含まれる全ての単語について、その単語と、当該単語が出現する本文における参照単位の位置の対を生成する。次に、分類処理機構13が、それぞれの単語について、その単語が表れる全ての参照単位に位置の組である（単語・参照単位の位置集合）対を得る。次に、検索インデックス生成部14により、全ての（単語・参照単位の位置集合）対から、各々の単語についての参照単位の位置集合を生成し、検索インデックス4を作成する。

【0047】これにより、文書内の全ての単語は、その直前のタグ（参照単位）の位置の情報と共に、検索インデックス4に登録されるので、検索対象とする単語から検索インデックス4を用いて検索することにより、文書中の単語の参照単位のタグ位置が直ちに検索でき、高速に参照単位の文書内容の部分が表示出力される。

【0048】図2は、第1の実施例の本文ファイルに対する検索インデックスの作成処理を示す処理フローを示すPAD（Problem Analysis Diagram）図である。また、図3は、図2に示す処理フローにより検索インデックスを作成する場合の作成プロセスの要部を具体的に説明する図である。図2および図3を参照して、全文登録語検索のための検索インデックスの作成処理を説明する。

【0049】まず、図2を参照して、本文ファイルに対する検索インデックスの作成処理の処理フローを概要を説明する。処理を開始すると、処理ブロック21において、本文ファイル31をタグの位置で分割し、分割した部分を参照単位とし、そのタグ位置を一時記憶する処理を行う。次に、繰り返し処理の制御ブロック22の処理を行う。この制御ブロック22の処理では、タグで分割された全てのタグ位置とその本文部分の対に対して、次の処理ブロック23および処理ブロック24の処理を繰り返し行う処理制御を行う。

【0050】この制御ブロック22の制御下の繰り返し処理では、まず、処理ブロック23において、処理対象の参照単位のタグ位置をAファイルに書き出す。次に、処理ブロック24において、処理対象の参照単位の本文の単語を、Aファイルに先に書き出したタグ位置に続いて、順番にAファイルに書き出す。これにより、Aファイルには、1つの参照単位について、タグ位置に続いて、その本文中の単語が連続して書き出される。このような処理を全ての参照単位について、制御ブロック22の処理制御により、繰り返し行う。このため、Aファイルには、図3に示すように、本文ファイル31から各々の参照単位について、まず、タグ位置が書き出されて、続いて当該タグ位置に対応する参照単位の本文の中の単語が順次書き出される。この結果、Aファイル32の内容は、タグ位置とそれに続く単語の組32aが、参照

単位の数だけ続くデータが得られる。

【0051】このようにして、Aファイルが作成されると、次に、繰り返し処理の制御ブロック25の処理を行う。この制御ブロック25の処理では、Aファイルに含まれる全ての単語に対して、次の処理ブロック26および処理ブロック27の処理を繰り返し行う処理制御を行う。

【0052】この制御ブロック25の制御下の繰り返し処理では、まず、処理ブロック26において、単語をキーとして、当該単語に対応するタグ位置を値とする対を作成する。続いて、処理ブロック27において、同じキー（単語）を持つ値（タグ位置）の対を集めて、キーと値の集合から構成されるリストを作成し、これをBファイルに書き出す。これにより、Bファイルには、1つの単語について、その単語が出現する参照単位のタグ位置のリストが得られる。このような処理を全てのAファイルの単語について、制御ブロック25による処理制御により、繰り返し行う。

【0053】この結果、図3に示すように、Bファイル33には、本文ファイルの各タグに区切られる参照単位の全ての単語について、当該単語がその出現する各々の参照単位に対応するタグ位置のリストが得られる。図3に示すBファイル33の例で説明すると、第1番目の単語1および第2番目の単語2に対応して、それぞれに『（単語1, 0, …）』および『（単語2, 0, 100, …）』のリストデータが得られている。つまり、これらのリストデータは、それぞれに『単語1が出現する参照単位のタグ位置がアドレス“0”, …であること』および『単語2が出現する参照単位のタグ位置がアドレス“0”, アドレス“100”, …であること』を意味している。

【0054】次に、処理ブロック28の処理を行い、Bファイルの内容に基づいて、単語からタグ位置の集合を検索できる検索インデックスを作成し、ここでの処理を終了する。これにより、各々の単語に対する検索インデックスが作成されると、その検索インデックスを用いることにより、検索対象の単語から直ちに、その単語が出現する参照単位のタグ位置の集合が得られる。したがって、検索単語から得られたタグ位置の集合に従って、当該タグ位置の集合からそれぞれの参照単位を表示できる。

【0055】以上に説明した全文登録語検索装置の第1の実施例においては、本文ファイルの参照単位をタグにより区分し、その位置を指示するタグ位置と、その中に含まれる単語を求めて記録する場合（Aファイル）、最初に参照単位の開始を指示するタグ位置を置き、続いて、その参照単位に属する単語を書くファイル形式をとっているが、各々の単語と参照単位（タグのタイプ）の間の関係を明確にして、同じ種類のタグの参照単位の中の単語を他と区別するため、各々の単語のデータにタグの種

類を示すフィールドを設けるようにしてもよい。これにより、同じ種類のタグの参照単位を検索単位として扱える。このような例を第2の実施例として説明する。

【0056】図4は、第2の実施例の全文登録語検索装置の本文ファイルに対する検索インデックスの作成処理を示す処理フローを示すPAD図である。また、図5は、図4に示す処理フローによる検索インデックスの作成プロセスの要部を具体的に説明する図である。図4および図5を参照して、第2の実施例の全文登録語検索のための検索インデックスの作成処理を説明する。

【0057】第2の実施例においては、単語から検索された結果のタグ位置により、表示する参照単位の区切りのタグの種類が直ちに判別できるように、参照単位ごとにその本文の検索対象となる単語の前にタグの種類を示すフィールドの文字列を付加している。これは、例えば、単語位置抽出部において、抽出した単語の前にタグ種別を示すフィールドの文字列を付加する処理を追加するように変形することにより、容易に対応できる。この種のタグの種類を示すフィールドを用いる場合の一例として、例えば、本文ファイルの参照単位の内容が、故障の個々の内容を示している場合に、故障の症状、原因、対処の3つのフィールドを1レコード中に設けておき、そのようなレコードが繰返し現われるような文書において、その単語の検索範囲を、症状を示すフィールドだけに限定する場合などに利用できる。この場合、タグの種類を示すフィールドの文字列として、症状、原因、対処の3つの種類を示す文字列を付加する。

【0058】図4を参照して、本文ファイルに対する検索インデックスの作成処理の処理フローを概要を説明する。処理を開始すると、処理ブロック41において、本文ファイルをタグの位置で分割し、分割した部分を参照単位とし、そのタグ位置を一時記憶する処理を行う。次に、繰返し処理の制御ブロック42の処理を行う。この制御ブロック42の処理では、タグで分割された全てのタグ位置とその本文部分の対に対して、次の処理ブロック43および処理ブロック44の処理を繰返し行う処理制御を行う。

【0059】この制御ブロック42の制御下の繰返し処理では、まず、処理ブロック43において、処理対象の参照単位のタグ位置をCファイルに書き出す。次に、処理ブロック44において、処理対象の参照単位の本文の各々の単語に対して、本文の単語の前にタグの種類を示すフィールドの文字列を付加したものを、新たな単語として、Cファイルに先に書き出したタグ位置に続いて、順番にCファイルに書き出す。これにより、Cファイルには、1つの参照単位について、タグ位置に続いて、タグの種類を示すフィールドの文字列を付加した本文中の単語が連続して書き出される。

【0060】このような処理を全ての参照単位について、制御ブロック42の処理制御により、繰返し行

う。この結果、図5に示すように、Cファイル52には、本文ファイル51から各々の参照単位について、まず、タグ位置が書き出されて、続いて当該タグ位置に対応する参照単位のタグの種類を示すフィールドの文字列(fld1など)を前に付加した本文の中の単語(単語1, 単語2など)が順次書き出される。この結果、Cファイル52の内容として、タグ位置とそれに続くタグの種類を示すフィールドの文字列を付加した単語の組52aが、参照単位の数だけ続くデータが得られる。

10 【0061】このようにして、Cファイルが作成されると、次に、繰返し処理の制御ブロック45の処理を行う。制御ブロック45の処理では、Cファイルに含まれる全ての単語に対して、次の処理ブロック46および処理ブロック47の処理を繰返し行う処理制御を行う。

【0062】この制御ブロック45の制御下の繰返し処理では、まず、処理ブロック46において、単語をキーとして、当該単語に対応するタグ位置を値とする対を作成する。続いて、処理ブロック47において、同じキー(単語)を持つ値(タグ位置)の対を集めて、キーと値の集合から構成されるリストを作成し、これをDファイルに書き出す。これにより、Dファイルには、前にタグの種類を示すフィールドの文字列を付加した1つの単語について、その単語が出現する参照単位のタグ位置のリストが得られる。このような処理を全てのCファイルに書き出された単語について、制御ブロック45による処理制御により、繰返し行う。

【0063】この結果、図5に示すように、Dファイル53には、本文ファイル51の全ての単語について、前にタグの種類を示すフィールドの文字列が付加された単語毎に、当該単語がその出現する個々の参照単位に対応して、そのタグ位置のリストが得られる。図5に示す例で説明すると、第1番目の単語1および第2番目の単語2に対しては、それぞれ『(fld1-単語1, 0, ...)』および『(fld1-単語2, 100, ...)』のリストデータが得られている。つまり、これらのリストデータは、『タグの種類が“fld1”である参照単位で単語1が出現するタグ位置が、アドレス“0”, ...であること』および『タグの種類が“fld1”である参照単位で単語2が出現するタグ位置がアドレス“0”, アドレス“100”, ...であること』をそれぞれ意味している。

【0064】次に、処理ブロック48の処理を行い、Dファイルの内容に基づいて、単語からタグ位置(タグフィールド名)の集合を検索できる検索インデックスを作成して、ここでの処理を終了する。これにより、各々の単語に対する検索インデックスが作成されると、その検索インデックスを用いることにより、タグの種類と検索対象の単語を指定することにより、タグの種類に応じて異なる参照単位についての検索対象の単語から直ちに、その単語が出現する参照単位のタグ位置の集合が得られる。したがって、検索単語から得られたタグ位置の集合

に従って、当該タグ位置からそれぞれの参照単位を表示できる。

【0065】以上に説明した全文登録語検索装置の第2の実施例においては、本文ファイルの参照単位をタグにより区分し、タグの種類に応じて、その位置を指示するタグ位置と、その中に含まれる単語を求めて記録する場合（Cファイル）、最初に参照単位の開始を指示するタグ位置を置き、続いて、その参照単位に属する単語に対しては、タグの種類を示すフィールドの文字列を付加して書くファイル形式をとっている。これにより、タグの種類に応じて、それぞれの参照単位の中の検索対象の単語の検索を、タグで区切る参照単位毎に高速に行うことができる。

【0066】また、タグを検索結果の参照単位の区切りとする全文登録語検索を行う場合において、本文の検索対象とする単語の中で、属性と値の対の記述が存在するものについては、その属性と値の対を検索対象の登録語として、登録しておくことにより、これらの単語の属性と値の対を検索対象として、前述の実施例と同様に、高速にタグの区切りを参照単位とする検索を行うことができる。このような実施例を、第3の実施例として説明する。

【0067】第3の実施例の全文登録語検索装置においては、タグを検索結果の参照単位の区切りとする全文登録語検索を行う場合、本文ファイルの本文の単語に、属性と値の対の記述のあるものについては、その対の記述を単語として登録する。具体的に説明すると、ここでの属性と値の対の例としては、本文中に、例えば、{売上高=100000}などのように、特別な形式により単語が示され、その単語が他と区別されている場合などがある。

【0068】図6は、第3の実施例の全文登録語検索装置の本文ファイルに対する検索インデックスの作成処理を示す処理フローを示すPAD図である。また、図7は、図6に示す処理フローにより検索インデックスを作成する作成プロセスの要部を具体的に説明する図である。図6および図7を参照して、第3の実施例の全文登録語検索のための検索インデックスの作成処理を説明する。

【0069】まず、図6を参照して、第3の実施例の本文ファイルに対する検索インデックスの作成処理の処理フローを概要を説明する。処理を開始すると、処理ブロック61において、本文ファイルをタグの位置で分割し、分割した部分を参照単位とし、そのタグ位置を一時記憶する処理を行う。次に、繰り返し処理の制御ブロック62の処理を行う。この制御ブロック62の処理では、タグで分割された全てのタグ位置とその本文部分（参照単位）の対に対して、次の処理ブロック63および処理ブロック64の処理を繰り返し行う処理制御を行う。

【0070】制御ブロック62の制御下の繰り返し処理では、まず、処理ブロック63において、処理対象の参照単位のタグ位置をEファイルに書き出す。次に、処理ブロック64において、処理対象の参照単位の本文の単語を、Eファイルに先に書き出したタグ位置に続いて、順番にEファイルに書き出す。ただし、属性と値の対の記述のあるものについては、その対も単語として、順番にEファイルに書き出す。これにより、Eファイルには、1つの参照単位について、当該タグ位置に続いて、その本文中の単語と、あれば属性と値の対とが連続して書き出される。このような処理を全ての参照単位について、制御ブロック62の処理制御により、繰り返し行う。

【0071】この結果、図7に示すように、Eファイル72には、本文ファイル71から各々の参照単位について、まず、タグ位置が書き出されて、続いて当該タグ位置に対応する参照単位の本文の中の単語と、属性と値の対とが順次書き出される。このため、前述の場合と同様に、Eファイル72の内容は、タグ位置とそれに続く単語と、属性と値の対との組が、参照単位の数だけ続くデータが得られる。

【0072】このようにして、Eファイルが作成されると、次に、繰り返し処理の制御ブロック65の処理を行う。制御ブロック65の処理では、Eファイルに含まれる全ての単語（属性と値の対を含む）に対して、次の処理ブロック66および処理ブロック67の処理を繰り返し行う処理制御を行う。

【0073】この制御ブロック65の制御下の繰り返し処理では、まず、処理ブロック65において、単語をキーとして、当該単語に対応するタグ位置を値とする対を作成する。続いて、処理ブロック67において、同じキー（単語）を持つ値（タグ位置）の対を集めて、キーと値の集合から構成されるリストを作成し、これをFファイルに書き出す。これにより、Fファイルには、1つの単語について、その単語が出現する参照単位のタグ位置のリストが得られる。このような処理を全てのEファイルの単語（属性と値の対を含む）について、制御ブロック65による処理制御により、繰り返し行う。

【0074】この結果、Fファイル73には、図7に示すように、本文ファイル71における全ての単語（属性と値の対を含む）について、当該単語がその出現する各々の参照単位に対応するタグ位置のリストが得られる。ここで、図7に示すFファイル73の例で説明すると、第1番目の単語1および第2番目の単語2に対応して、それぞれに『（単語1，0，…）』および『（単語2，0，100，220，…）』のリストデータが得られている。つまり、これらのリストデータは、それぞれに『単語1が出現する参照単位のタグ位置がアドレス“0”，…であること』および『単語2が出現する参照単位のタグ位置がアドレス“0”，アドレス“10

0", アドレス "220", ...であること』を意味している。また、ここでは、属性と値の対についても、単語の場合と同様に『({××事業規模=1000}, 0, ...)』のリストデータが得られており、このリストデータは、『属性と値の対 {××事業規模=1000} が出現する参照単位のタグ位置がアドレス "0" ...であること』を意味している。

【0075】次に、処理ブロック68の処理を行い、Fファイルの内容に基づいて、単語からタグ位置の集合を検索できる検索インデックスを作成し、ここでの処理を終了する。これにより、各々の単語に対する検索インデックスが作成されると、その検索インデックスを用いることによって、検索対象として単語を指示することにより、前述の場合と同様に、検索対象の単語が出現する参照単位のタグ位置の集合が得られる。また、例えば、検索対象として属性と値の対を指示することにより、直ちに、検索対象の属性と値の対が出現する参照単位のタグ位置の集合が得られる。したがって、検索単語から得られたタグ位置の集合に従って、当該タグ位置の集合からそれぞれの参照単位を表示できる。

【0076】次に、前述した第2の実施例による単語の前にタグの種類を示すフィールドの文字列を付加し、更に、第3の実施例による属性と値の対のあるものについては、その属性と値の対を単語として登録する場合の変形を組合せるようにしても良い。このような実施例について、第4の実施例として説明する。

【0077】つまり、第4の実施例による全文登録語検索装置は、タグを検索結果の参照単位の区切りとする全文登録語検索装置のうち、属性と値の対の記述のあるものについては、その対を単語として登録し、その際、単語として登録する属性と値の対の前にタグの種類を示すフィールドの文字列を付加して登録する。

【0078】図8は、第4の実施例の全文登録語検索装置の本文ファイルに対する検索インデックスの作成処理を示す処理フローを示すPAD図である。また、図9は、図8に示す処理フローによる検索インデックスの作成プロセスの要部を具体的に説明する図である。図8および図9を参照して、第4の実施例の全文登録語検索のための検索インデックスの作成処理を説明する。

【0079】第4の実施例においては、単語または属性と値の対の指定により、検索された結果のタグ位置により、表示する参照単位の区切りのタグの種類が直ちに判別できるように、参照単位ごとにその本文の検索対象となる単語(属性と値の対を含む)の前にタグの種類を示すフィールドの文字列を付加する。これも、前述したように、例えば、単語位置抽出部において、抽出した単語の前にタグ種別を示すフィールドの文字列を付加する処理を追加するように変形することにより、容易に対応できる。

【0080】図8を参照して、本文ファイルに対する検

索インデックスの作成処理の処理フローを概要を説明する。処理を開始すると、処理ブロック81において、本文ファイルをタグの位置で分割し、分割した部分を参照単位とし、そのタグ位置を一時記憶する処理を行う。次に、繰り返し処理の制御ブロック82の処理を行う。この制御ブロック82の処理では、タグで分割された全てのタグ位置とその本文部分の対に対して、次の処理ブロック83および処理ブロック84の処理を繰り返し行う処理制御を行う。

10 【0081】この制御ブロック82の制御下の繰り返し処理では、まず、処理ブロック83において、処理対象の参照単位のタグ位置をGファイルに書き出す。次に、処理ブロック84において、処理対象の参照単位の本文の各々の単語に対して、本文の単語の前にタグの種類を示すフィールドの文字列を付加したものを、新たな単語として、Gファイルに先書き出したタグ位置に続いて、順番にGファイルに書き出す。ただし、この場合、属性と値の対の記述のあるものについても、その対を単語として、順番にGファイルに書き出す。これにより、G

20 ファイルには、1つの参照単位について、タグ位置に続いて、タグの種類を示すフィールドの文字列を前に付加した本文中の単語または属性と値の対が連続して書き出される。

【0082】このような処理を全ての参照単位について、制御ブロック82の処理制御により、繰り返し行う。この結果、図9に示すように、Gファイル92には、本文ファイル91から各々の参照単位について、まず、タグ位置が書き出されて、続いて当該タグ位置に対応する参照単位のタグの種類を示すフィールドの文字列を前に付加した本文の中の単語が順次書き出され、または、タグの種類を示すフィールドの文字列を付加した本文の中の属性と値の対が書き出される。この結果、G

30 ファイル92の内容として、タグ位置とそれに続くタグの種類を示すフィールドの文字列を付加した単語または属性と値の対の組が、参照単位の数だけ続くデータが得られる。

【0083】このようにして、Gファイルが作成されると、次に、繰り返し処理の制御ブロック85の処理を行う。制御ブロック85の処理では、Gファイルに含まれる全ての単語(属性と値の対を含む)に対して、次の処理

40 ブロック86および処理ブロック87の処理を繰り返し行う処理制御を行う。

【0084】この制御ブロック85の制御下の繰り返し処理では、まず、処理ブロック86において、単語(属性と値の対を含む)をキーとして、当該単語に対応するタグ位置を値とする対を作成する。続いて、処理ブロック87において、同じキー(単語、属性と値の対)を持つ値(タグ位置)の対を集めて、キーと値の集合から構成されるリストを作成し、これをHファイルに書き出す。これにより、Hファイルには、タグの種類を示すフ

フィールドの文字列を前に付加した1つの単語について、その単語が出現する参照単位のタグ位置のリストが得られる。このような処理を全てのGファイルに書き出されている単語（属性と値の対を含む）について、制御ブロック85による処理制御により、繰り返し行う。

【0085】この結果、Hファイル93には、図9に示すように、本文ファイル91の全ての単語について、タグの種類を示すフィールドの文字列を前に付加した1つの単語毎に、当該単語がその出現する個々の参照単位に対応して、そのタグ位置のリストが得られる。図9に示すHファイル93の例で説明すると、第1番目の単語1および第2番目の単語2に対しては、それぞれ『(fld1-単語1, 0, …)』および『(fld1-単語2, 0, 220, …)』のリストデータが得られている。つまり、これらのリストデータは、『タグの種類が“fld1”である参照単位で単語1が出現するタグ位置が、アドレス“0”, …であること』および『タグの種類が“fld1”である参照単位で単語2が出現するタグ位置がアドレス“0”, アドレス“220”, …であること』をそれぞれ意味している。また、属性と値の対についても、単語の場合と同様な形式で『(fld1- {××事業規模=1000}, 0, …)』のリストデータが得られており、このリストデータは、『タグの種類が“fld1”である参照単位で、属性と値の対 {××事業規模=1000} が出現するタグ位置がアドレス“0” …であること』を意味している。

【0086】次に、処理ブロック88の処理を行い、Hファイルの内容に基づいて、単語からタグ位置（タグフィールド名）の集合を検索できる検索インデックスを作成して、ここでの処理を終了する。これにより、各々の単語に対する検索インデックスが作成されると、その検索インデックスを用いることにより、タグの種類と検索対象の単語として属性と値の組を指定することにより、タグの種類に応じて異なる参照単位についての検索対象の属性の値の組から直ちに、その属性と値の組が出現する参照単位のタグ位置の集合が得られる。したがって、検索単語から得られたタグ位置の集合に従って、当該タグ位置からそれぞれの参照単位を表示できる。

【0087】以上、本発明の実施例について説明したが、本発明は、上述した実施例に限定されるものではない。ここでは、全ての単語を抽出して登録語とする例について示しているが、助詞、助動詞、接続詞など、検索対象を特定して検索語として登録するようにしてもよいことは明らかである。その際、検索対象の特徴を必ずしも示していないような単語を登録しないことも可能である。

【0088】

【発明の効果】以上、説明したように、本発明の全文登録語検索装置によれば、本文中の検索対象とする単語に

ついて、例えば、本文中の全ての単語について、その単語が出現する位置の直前にあるタグの位置を全て集めて、検索インデックスを作成して検索装置を構成する。そして、検索対象の単語を入力とし、その単語が出現する位置の直前にあるタグの位置の集合を検索結果として出力する。これにより、タグを有する文書において、タグを検索結果の参照単位の区切りとして、文書内の全ての単語が、その直前のタグの位置と共に、検索インデックスとして保存されるので、文書内の全ての単語に対し、当該単語を含むタグで区切られた参照単位を即座に検索することが可能となる。

【図面の簡単な説明】

【図1】 図1は本発明の実施例の全文登録語検索装置の装置構成の要部を示すブロック図、

【図2】 図2は第1の実施例の本文ファイルに対する検索インデックスの作成処理を示す処理フローを示すPAD (Problem Analysis Diagram) 図、

【図3】 図3は図2に示す処理フローにより検索インデックスを作成する場合の作成プロセスの要部を具体的に説明する図、

【図4】 図4は第2の実施例の全文登録語検索装置の本文ファイルに対する検索インデックスの作成処理を示す処理フローを示すPAD図、

【図5】 図5は図4に示す処理フローによる検索インデックスの作成プロセスの要部を具体的に説明する図、

【図6】 図6は第3の実施例の全文登録語検索装置の本文ファイルに対する検索インデックスの作成処理を示す処理フローを示すPAD図、

【図7】 図7は図6に示す処理フローによる検索インデックスの作成プロセスの要部を具体的に説明する図、

【図8】 図8は第4の実施例の全文登録語検索装置の本文ファイルに対する検索インデックスの作成処理を示す処理フローを示すPAD図、

【図9】 図9は図8に示す処理フローによる検索インデックスの作成プロセスの要部を具体的に説明する図、

【図10】 図10はSGMLの文書型定義(DTD)の一例を示す図、

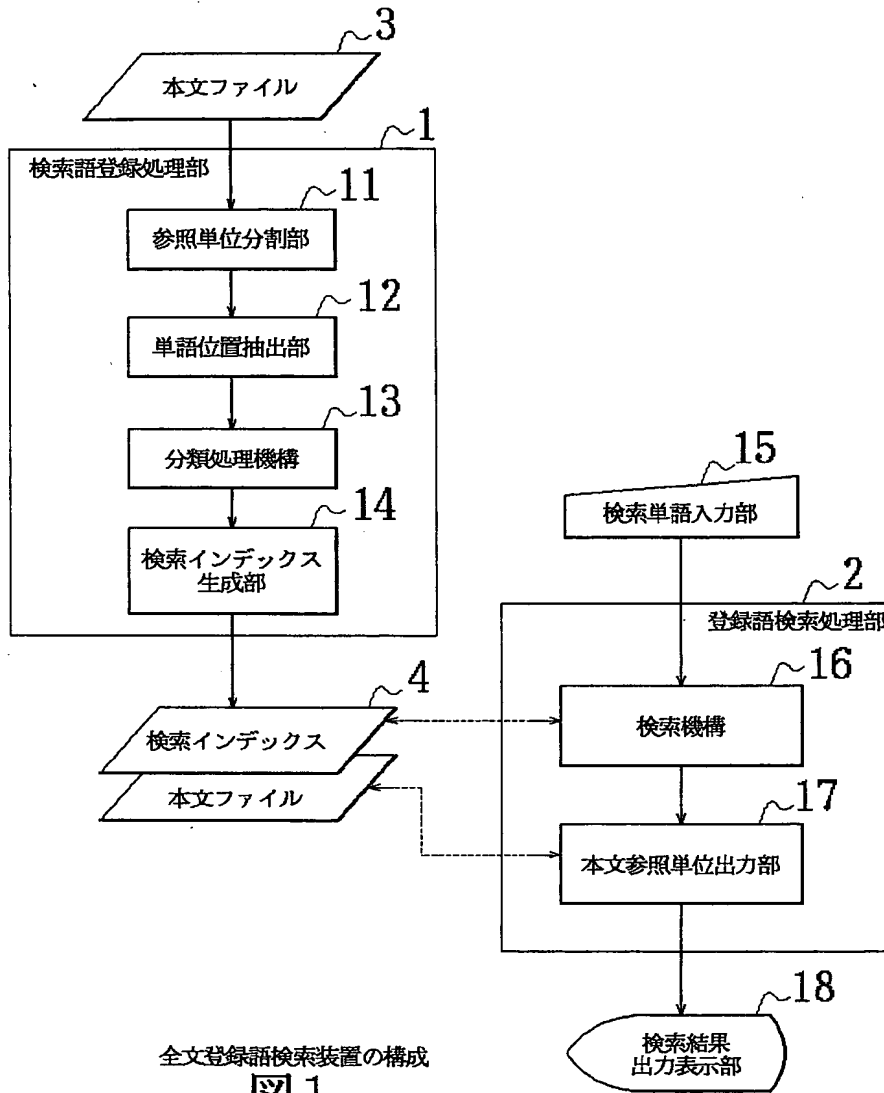
【図11】 図11はタグ付き文書としてのタグが省略されたSGML文書の一例を説明する図、

【図12】 図12は省略されたタグが復元されたSGML文書の一例を説明する図である。

【符号の説明】

1…検索語登録処理部、2…登録語検索処理部、3…タグを有する文書の本文ファイル、4…本文ファイルに付加された検索インデックス、11…参照単位分割部、12…単語位置抽出部、13…分類処理機構、14…検索インデックス生成部、15…検索単語入力部、16…検索機構、17…本文参照単位出力部、18…検索結果出力表示部。

【図1】



全文登録語検索装置の構成

図1

【図10】

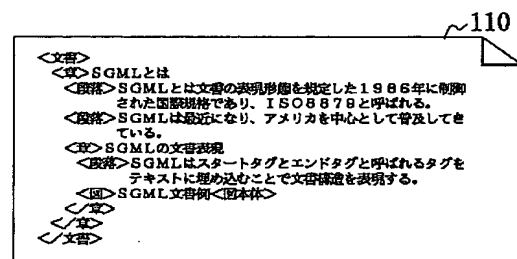
100

<!ELEMENT	文章	--	章*	>
<!ELEMENT	章	--	題、(段落 図)*、(章*)	>
<!ELEMENT	題	OO	(#PCDATA)	>
<!ELEMENT	段落	—O	(#PCDATA)	>
<!ELEMENT	図	—O	(図、図本体)	>
<!ELEMENT	図本体	—O	EMPTY	>

SGMLの文書型定義(DTD)の一例

図10

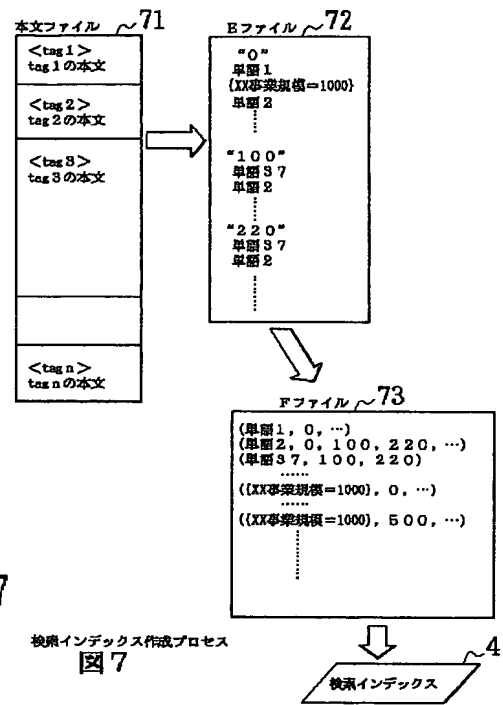
【図11】



SGMLによる文書の例

図11

【図 7】



本文ファイル 91

Gファイル 92

Hファイル 93

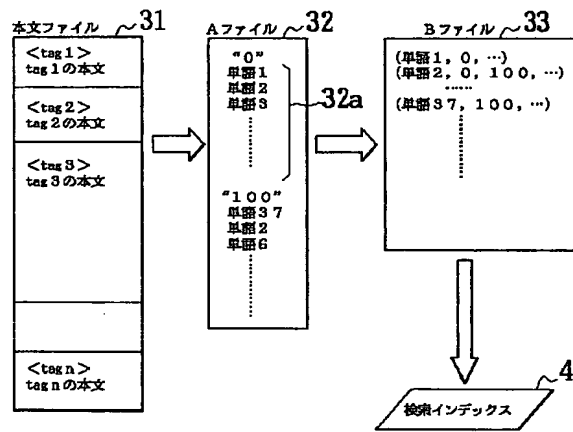
検索インデックス作成プロセス

図 9

[illegible]

— 13 —

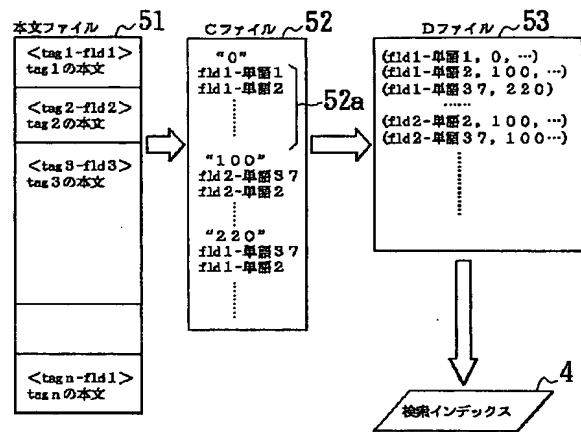
【図3】



検索インデックス作成プロセス

図3

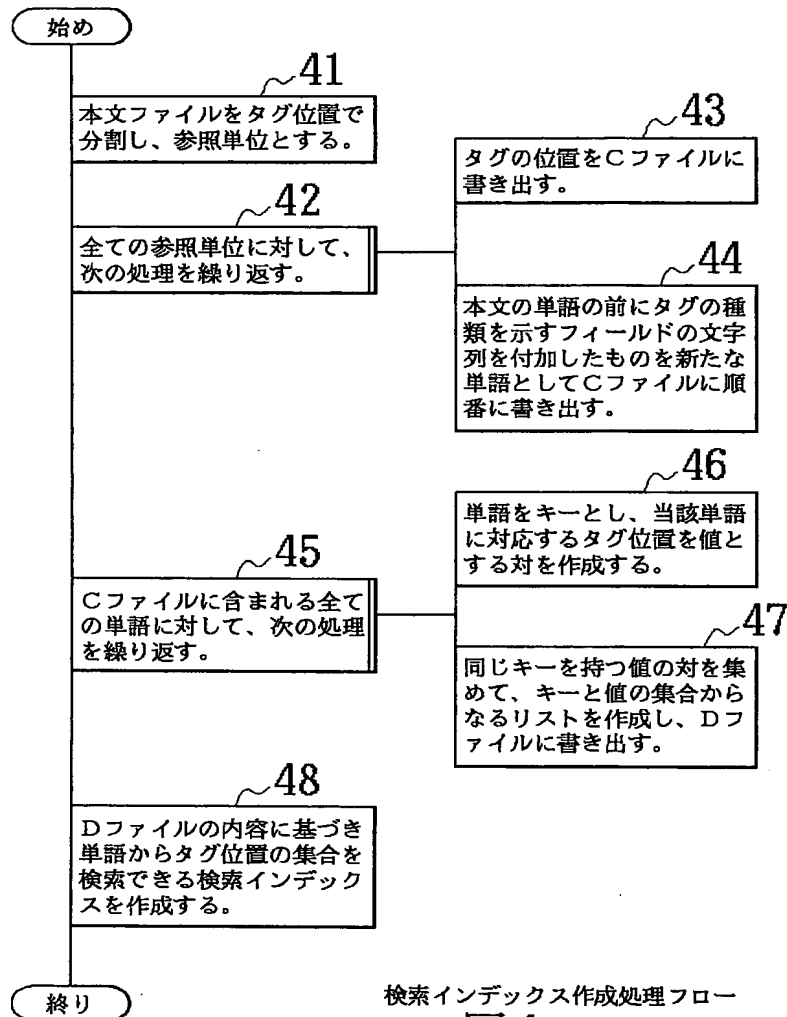
【図5】



検索インデックス作成プロセス

図5

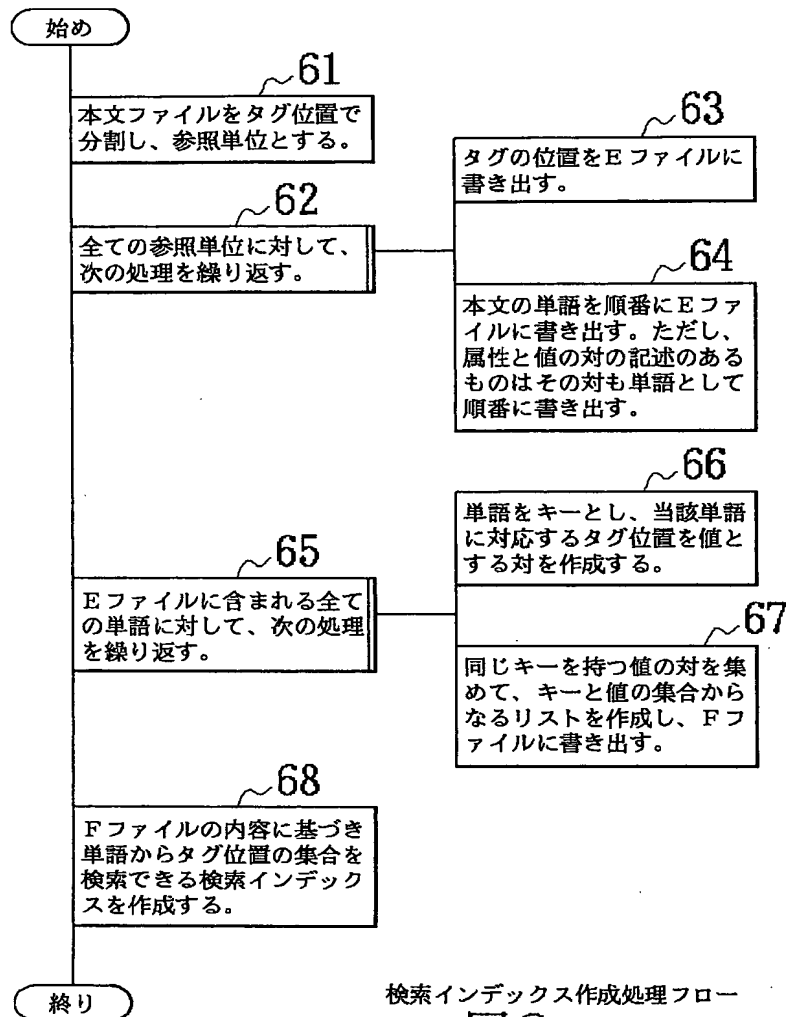
【図 4】



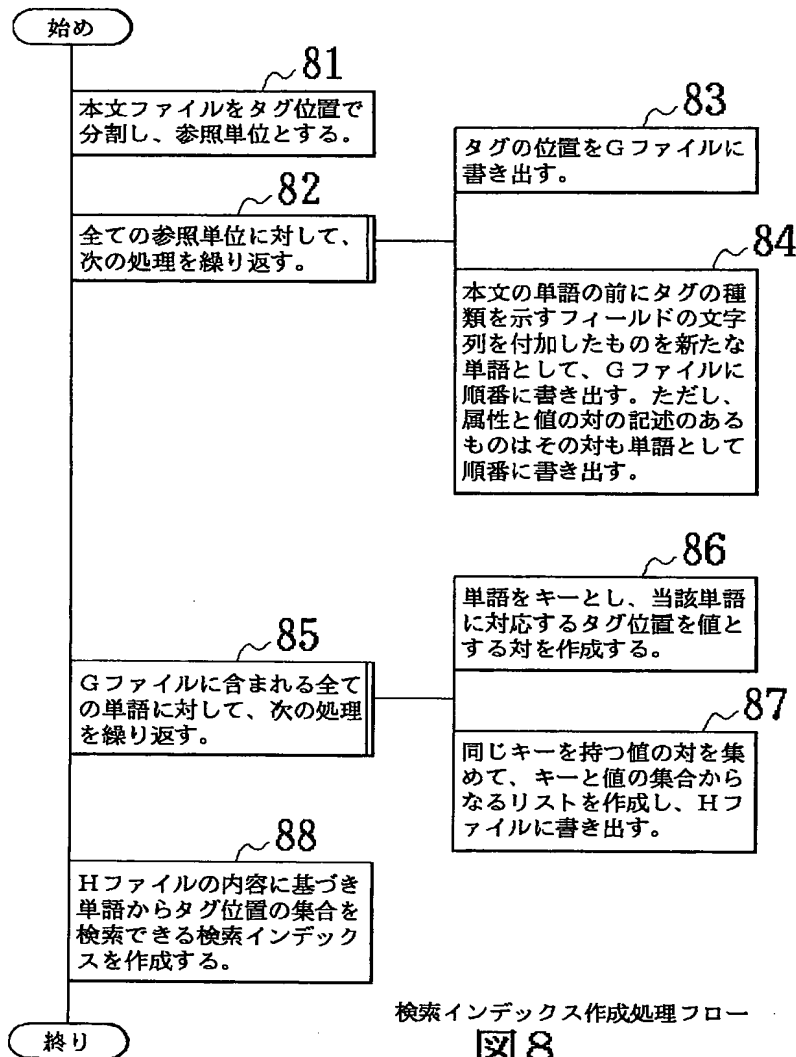
検索インデックス作成処理フロー

図 4

【図6】



【図 8】



**This Page is Inserted by IFW Indexing and Scanning
Operations and is not part of the Official Record**

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☐ BLACK BORDERS
- ☐ IMAGE CUT OFF AT TOP, BOTTOM OR SIDES
- ☒ FADED TEXT OR DRAWING
- ☒ BLURRED OR ILLEGIBLE TEXT OR DRAWING
- ☐ SKEWED/SLANTED IMAGES
- ☐ COLOR OR BLACK AND WHITE PHOTOGRAPHS
- ☐ GRAY SCALE DOCUMENTS
- ☐ LINES OR MARKS ON ORIGINAL DOCUMENT
- ☒ REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY
- ☐ OTHER: _____

IMAGES ARE BEST AVAILABLE COPY.

As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.